# Variational Treatment of Probabilistic Directed Graphical Models

Satish Kumar

July 27, 2018

LEAP Lab, IISC Bengaluru

## Table of contents

# Generative modeling

Figure 1: Some generative models perform density estimation. These models take a training set of examples drawn from an unknown data-generating distribution $p_{\text{data}}$ and return an estimate of that distribution. The estimate $p_{\text{model}}$ can be evaluated for a particular value of $\boldsymbol{x}$ to obtain an estimate $p_{\text{model}}(\boldsymbol{x})$ of the true density $p_{\text{model}}(\boldsymbol{x})$. This figure illustrates the process for a collection of samples of one-dimensional data and a Gaussian model.

# An Ideal Generative model



Training examples → Model samples

# Why study Generative modeling?

- Training and sampling from generative models is an excellent test of our ability to represent and manipulate high-dimensional probability distributions.
- Generative models can be trained with missing data and can provide predictions on inputs that are missing data.
- It enable machine learning to work with multi-modal outputs
- Realistic generation of samples from some distribution.

# Evidence Lower Bound(ELBO)

Consider a probabilistic model:Observed variables X and latent variables Z.Our goal is to minimize the likelihood function given by

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$

- We introduce a distribution q(z) defined over latent variables,and for any choice of q(z) the following decomposition holds

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \mathrm{KL}(q\|p)$$
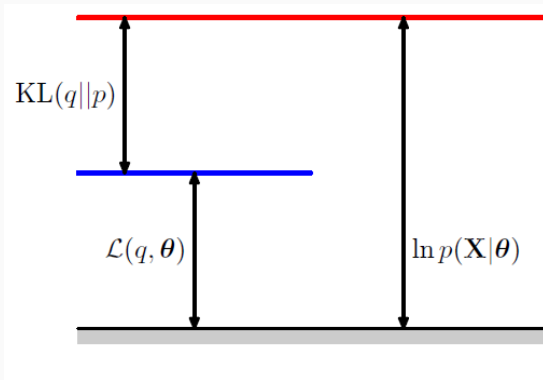
where we have defined

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$

$$\mathrm{KL}(q\|p) = -\sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}.$$

- using the product rule of probability

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) = \ln p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}) + \ln p(\mathbf{X}|\boldsymbol{\theta})$$
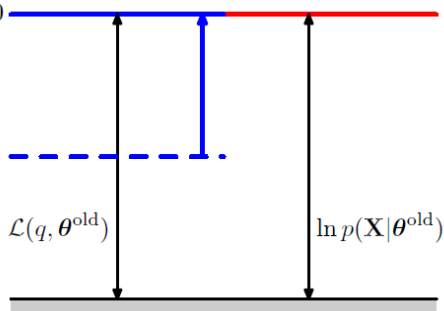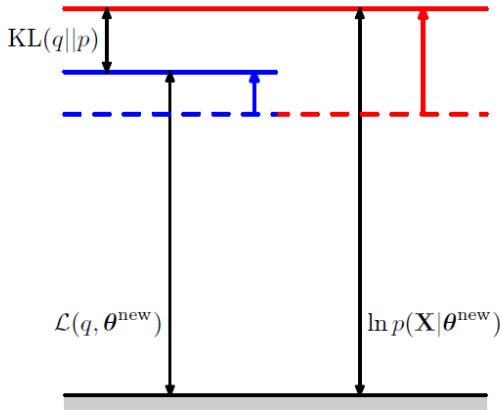
Illustration of the E step of the EM algorithm. The $q$ distribution is set equal to the posterior distribution for the current parameter values $\theta^{\mathrm{old}}$, causing the lower bound to move up to the same value as the log likelihood function, with the KL divergence vanishing.
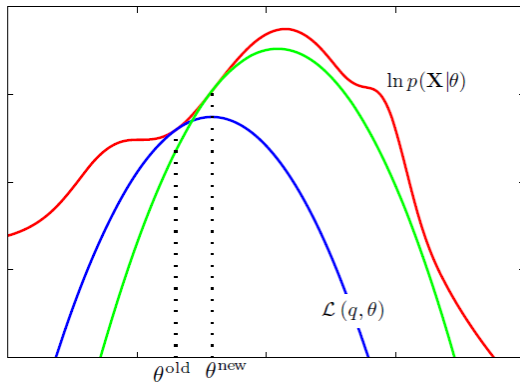
$$\mathrm{KL}(q\|p) = 0$$

$$\mathcal{L}(q, \boldsymbol{\theta}^{\mathrm{old}})$$

$$\ln p(\mathbf{X}|\boldsymbol{\theta}^{\mathrm{old}})$$

Illustration of the M step of the EM algorithm. The distribution $q(\mathbf{Z})$ is held fixed and the lower bound $\mathcal{L}(q, \boldsymbol{\theta})$ is maximized with respect to the parameter vector $\boldsymbol{\theta}$ to give a revised value $\boldsymbol{\theta}^{\text{new}}$. Because the KL divergence is nonnegative, this causes the log likelihood $\ln p(\mathbf{X}|\boldsymbol{\theta})$ to increase by at least as much as the lower bound does.



$\text{KL}(q||p)$

$\mathcal{L}(q, \boldsymbol{\theta}^{\text{new}})$

$\ln p(\mathbf{X}|\boldsymbol{\theta}^{\text{new}})$

The EM algorithm involves alternately computing a lower bound on the log likelihood for the current parameter values and then maximizing this bound to obtain the new parameter values. See the text for a full discussion.

# Variational Inference in Deep Learning

## Some preliminaries

- We assume the observed variable $x$ is a random sample from an unknown underlying process where the distribution $p^*(x)$ is unknown
- we approximate this process with a chosen model $p_\theta(x)$ with parameter $\theta$

$$x \sim p_\theta(x)$$

- Learnig is the process of searching the $\theta$ such that, for any observed variable $x$

$$p_\theta(x) \approx p^*(x)$$

- we wish $p_\theta(x)$ to be sufficiently flexible to be able to adapt to the data
- Often, such as in case of classification or regression problems, we are not interested in learning an unconditional model $p_\theta(x)$, but a conditional model $p_\theta(y|x)$
- that approximates the underlying conditional distribution $p^*(y|x)$

$$p_\theta(y|x) \approx p^*(y|x)$$

## Parameterizing Conditional distributions with Neural Nets

- We parameterize conditional distributions with neural networks
- In image classification, neural networks parameterize a categorical distribution $p_\theta(y|x)$ over a class label $y$, conditioned on an image $x$.

$$\mathbf{p} = \text{NeuralNet}(\mathbf{x})$$
$$p_\theta(y|\mathbf{x}) = \text{Categorical}(y; \mathbf{p})$$

- where for all $p_i \in \mathbf{p}$
- and the last operation of the neural net is a *softmax()* function such that $\sum_i p_i = 1$

In the case of a Gaussian MLP as encoder or decoder, we let the encoder or decoder be a multivariate Gaussian with a diagonal covariance structure:

$$\log p(\mathbf{x}|\mathbf{z}) = \log \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \sigma^2 \mathbf{I})$$
$$\text{where } \boldsymbol{\mu} = \mathbf{W}_4 \mathbf{h} + \mathbf{b}_4$$
$$\log \sigma^2 = \mathbf{W}_5 \mathbf{h} + \mathbf{b}_5$$
$$\mathbf{h} = \tanh(\mathbf{W}_3 \mathbf{z} + \mathbf{b}_3)$$

where $\{\mathbf{W}_3, \mathbf{W}_4, \mathbf{W}_5, \mathbf{b}_3, \mathbf{b}_4, \mathbf{b}_5\}$ are the weights and biases of the MLP and part of $\theta$ when used as decoder. Note that when this network is used as an encoder $q_{\phi}(\mathbf{z}|\mathbf{x})$, then $\mathbf{z}$ and $\mathbf{x}$ are swapped, and the weights and biases are variational parameters $\phi$.

## Learning in Fully Observed models with neural nets

- **Dataset**: We often collect a dataset $D$ consisting of $N \geq 1$ datapoints:

$$\mathcal{D} = \{x^{(1)}, x^{(2)}, ..., x^{(N)}\} \equiv \{x^{(i)}\}_{i=1}^N \equiv x^{(1:N)}$$

- Under the i.i.d. assumption, the probability of the datapoints given the parameters factorizes as a product of individual datapoint probabilities.

- **Maximum Likelihood and Minibatch SGD**

$$\log p_\theta(\mathcal{D}) = \sum_{x \in \mathcal{D}} \log p_\theta(x)$$

- Using automatic differentiation tools, we can efficiently compute gradients of this objective, and use such gradient to find the local optimum of the ML objective.

- We can opt either Stochastic Gradient or Batch Gradient Methods

## Learning and Inference in Deep latent variable models

- **Latent Variables, z** are variables that are part of the model, but which we dont observe, and are therefore **not part of the dataset**.
- In case of unconditional modeling of observed variable $x$, the directed graphical model would then represent a joint distribution $p_\theta(x,z)$ over both the observed variables $x$ and the latent variables $z$.
- The marginal distribution over the observed variables $p_\theta(x)$, is given by:

$$p_\theta(x) = \int p_\theta(x,z)dz$$

- This is also called the (single datapoint) **marginal likelihood** or the **model evidence**, when taking as a function of $\theta$.

## Deep Latent Variable Models(DLVMs)

- We use the term deep latent variable model (DLVM) to denote a latent variable model $p_\theta(\boldsymbol{x},\boldsymbol{z})$ whose distributions are parameterized by neural networks

- Such a model can be conditioned on some context, like $p_\theta(\boldsymbol{x},\boldsymbol{z}|\boldsymbol{y})$

- The simplest, and most common, graphical model with latent variables is one that is specified as factorization with the following structure:

$$p_\theta(\boldsymbol{x},\boldsymbol{z})=p_\theta(\boldsymbol{z})p_\theta(\boldsymbol{x}|\boldsymbol{z})$$

- The distribution $p(\boldsymbol{z})$ is often called the prior distribution over $\boldsymbol{z}$

## Example: DLVM for multivariate Bernoulli data

- For binary data $\mathbf{x}$, with a spherical Gaussian latent space, and a factorized Bernoulli observation model:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, \mathbf{I})$$
$$\mathbf{p} = \text{DecoderNeuralNet}_{\theta}(\mathbf{z})$$
$$\log p(\mathbf{x}|\mathbf{z}) = \sum_{j=1}^{D} \log p(x_j|\mathbf{z}) = \sum_{j=1}^{D} \log \text{Bernoulli}(x_j; p_j)$$
$$= \sum_{j=1}^{D} x_j \log p_j + (1 - x_j) \log(1 - p_j)$$

- where $\forall p_j \; \mathbf{p} : 0 \leq p_j \leq 1$ (e.g. implemented through a sigmoid nonlinearity as the last layer of the *DecoderNeuralNet*(.)), where D is the dimensionality of $\mathbf{x}$,
- *Bernoulli*(.; $p$) is the probability mass function (PMF) of the Bernoulli distribution.

## Intractabilities

- The main difficulty of maximum likelihood learning in DLVMs is that the marginal probability of data under the model is typically **intractable**.
- This is due to the integral in equation for computing the marginal likelihood (or modelevidence)
- Note that the joint distribution $p_\theta(x,z)$ is efficient to compute, and that the densities are related through the basic identity:

$$p_\theta(z|x) = \frac{p_\theta(x,z)}{p_\theta(x)}$$

- Since $p_\theta(x,z)$ is tractable to compute
- a tractable marginal likelihood $p_\theta(x)$ leads to a tractable posterior $p_\theta(z|x)$, and vice versa
- But both are intractable
- Approximate inference techniques allow us to approximate the posterior $p_\theta(z|x)$ and the marginal likelihood $p_\theta(x)$ in DLVMs.

# Variational Autoencoder

## Encoder or approximate posterior

- Let $p_\theta(\boldsymbol{x},\boldsymbol{z})$ be a latent-variable model with observed variables x and latent variables $\boldsymbol{z}$
- To turn a DLVMs intractable posterior inference and learning problems into tractable problems, we introduce a parametric inference model $q_\phi(\boldsymbol{z}|\boldsymbol{x})$
- This model is also called an encoder.
- With $\phi$ we indicate the parameters of this inference model, also called the variational parameters.
- We optimize the variational parameters such that:

$$q_\phi(\boldsymbol{z}|\boldsymbol{x}) \approx p_\theta(\boldsymbol{z}|\boldsymbol{x})$$

- As we will explain, this approximation to the posterior help us optimize the marginal likelihood.

- Like a DLVM, the inference model can be (almost) any directed graphical model:

$$q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) = q_{\boldsymbol{\phi}}(\mathbf{z}_1, ..., \mathbf{z}_M|\mathbf{x}) = \prod_{j=1}^{M} q_{\boldsymbol{\phi}}(\mathbf{z}_j|Pa(\mathbf{z}_j), \mathbf{x})$$

- similar to a DLVM, the distribution $q_{\phi}(\mathbf{z}|\mathbf{x})$ can be parameterized using deep neural networks

- In this case, the variational parameters $\boldsymbol{\phi}$ include the weights and biases of the neural network

$$(\boldsymbol{\mu}, \log \boldsymbol{\sigma}) = \text{EncoderNeuralNet}_{\boldsymbol{\phi}}(\mathbf{x})$$
$$q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}))$$

- we use a single encoder neural network to perform posterior inference over all of the datapoints in our dataset

## Evidence lower bound for VAE

- The **optimization objective** of the variational autoencoder, like in other variational methods, is the evidence lower bound
- Also called Variational lower bound
- For any choice of inference model $q_\phi(z|x)$, including the choice of variational parameters , we have:

$$
\begin{aligned}
\log p_\theta(\mathbf{x}) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log p_\theta(\mathbf{x}) \right] \\
&= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \left[ \frac{p_\theta(\mathbf{x}, \mathbf{z})}{p_\theta(\mathbf{z}|\mathbf{x})} \right] \right] \\
&= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \left[ \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} \right] \right] \\
&= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \left[ \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \right]}_{\substack{=\mathcal{L}_{\theta,\phi}(\mathbf{x}) \\ \text{(ELBO)}}} + \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \left[ \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} \right] \right]}_{=D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))}
\end{aligned}
$$

- KL divergence between $q_\phi(\mathbf{z}|\mathbf{x})$ and $p_\theta(\mathbf{z}|\mathbf{x})$, which is non-negative:
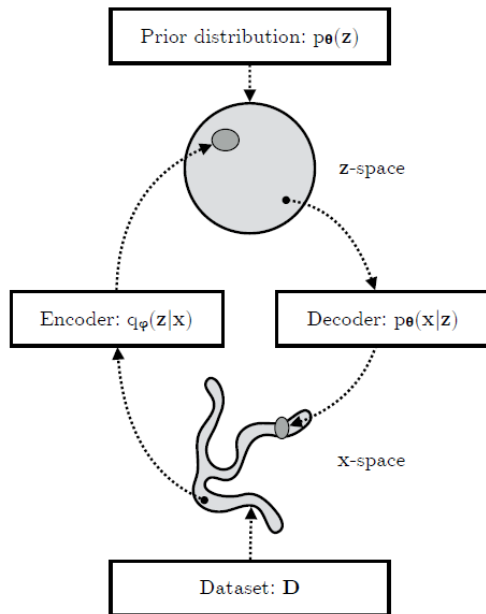$$D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) \geq 0$$

- and zero if, and only if, $q_\phi(\mathbf{z}|\mathbf{x})$ equals the true posterior distribution

- The first term in the variational lower bound, also called the evidence lower bound (ELBO):
$$\mathcal{L}_{\theta,\phi}(\mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}) \right]$$
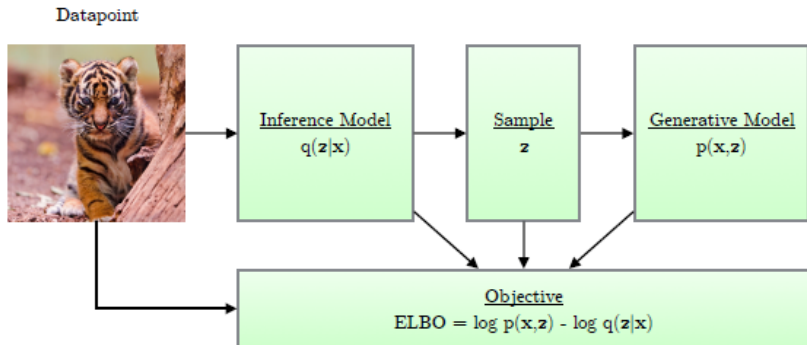
- Due to the non-negativity of the *KLDivergence*, the ELBO is a lower bound on the log-likelihood of the data.
$$\mathcal{L}_{\theta,\phi}(\mathbf{x}) = \log p_\theta(\mathbf{x}) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))$$
$$\leq \log p_\theta(\mathbf{x})$$

## A Double-Edged Sword

- Maximization of the ELBO $L_{\theta,\phi}(x)$ w.r.t. the parameters $\theta$ and $\phi$, will concurrently optimize the two things we care about:

- It will approximately maximize the marginal likelihood $p_\theta(x)$. This means that our generative model will become better.

- It will minimize the KL divergence of the approximation $q_\phi(z|x)$ from the true posterior $p_\theta(z|x)$, so $q_\phi(z|x)$ becomes better.

## Stochastic gradient-based optimization of the ELBO

- An important property of the ELBO, is that it allows joint optimization w.r.t. all parameters ($\phi$ and $\theta$) using SGD.
- We can start out with random initial values of $\phi$ and $\theta$, and stochastically optimize their values until convergence.
- Given a dataset with i.i.d. data, the ELBO objective is the sum (or average) of individual-datapoint ELBOs:

$$\mathcal{L}_{\theta,\phi}(\mathcal{D}) = \sum_{x \in \mathcal{D}} \mathcal{L}_{\theta,\phi}(x)$$

- Unbiased gradients of the ELBO w.r.t. the generative model parameters are simple to obtain:

$$\nabla_{\theta}\mathcal{L}_{\theta,\phi}(\mathbf{x}) = \nabla_{\theta}\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}\left[\log p_{\theta}(\mathbf{x},\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})\right]$$
$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}\left[\nabla_{\theta}(\log p_{\theta}(\mathbf{x},\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}))\right]$$
$$\simeq \nabla_{\theta}(\log p_{\theta}(\mathbf{x},\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}))$$
$$= \nabla_{\theta}(\log p_{\theta}(\mathbf{x},\mathbf{z}))$$

- Unbiased gradients w.r.t. the variational parameters are more difficult to obtain, since the ELBOs expectation is taken w.r.t. the distribution $q_{\phi}(\mathbf{z}|\mathbf{x})$, which is a function of $\phi$. I.e., in general:

$$\nabla_{\phi}\mathcal{L}_{\theta,\phi}(\mathbf{x}) = \nabla_{\phi}\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}\left[\log p_{\theta}(\mathbf{x},\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})\right]$$
$$\neq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}\left[\nabla_{\phi}(\log p_{\theta}(\mathbf{x},\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}))\right]$$

# Reparameterization trick

## Change of variables

- For continuous latent variables and a differentiable encoder and generative model, the ELBO can be straightforwardly differentiated w.r.t. both $\theta$ and $\phi$ through a **change of variables**, also called the **reparameterization trick**

- First, we express the random variable $\mathbf{z}$ as some differentiable (and invertible) transformation of another random variable $\epsilon$, given $\mathbf{z}$ and

$$\mathbf{z} = \mathbf{g}(\epsilon, \phi, \mathbf{x})$$

- where the distribution of random variable $\epsilon$ is independent of $\mathbf{x}$ or $\phi$

## Gradient of expectation under change of variable

- Given such a change of variable, expectations can be rewritten in terms of $\epsilon$

$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}\left[f(\mathbf{z})\right] = \mathbb{E}_{p(\epsilon)}\left[f(\mathbf{z})\right]$$

- where $\mathbf{z} = \mathbf{g}(\epsilon, \phi, \mathbf{x})$ and the expectation and gradient operators become commutative, and we can form a simple Monte Carlo estimator:
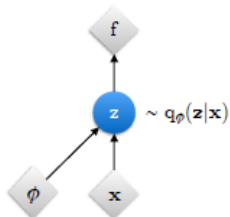
$$\begin{aligned}\nabla_{\phi}\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}\left[f(\mathbf{z})\right] &= \nabla_{\phi}\mathbb{E}_{p(\epsilon)}\left[f(\mathbf{z})\right] \\ &= \mathbb{E}_{p(\epsilon)}\left[\nabla_{\phi}f(\mathbf{z})\right] \\ &\simeq \nabla_{\phi}f(\mathbf{z})\end{aligned}$$

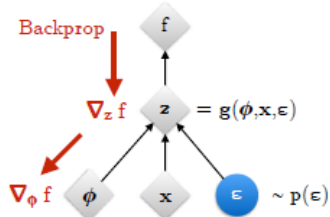- where in the last line, $\mathbf{z} = \mathbf{g}(\epsilon, \phi, \mathbf{x})$ with random noise sample

$$\epsilon \sim p(\epsilon)$$

# Illustration of the reparameterization trick

## Gradient of ELBO

- Under the reparameterization, we can replace an expectation w.r.t. $q_\phi(z|x)$ with one w.r.t. $p(\epsilon)$. The ELBO can be rewritten as:

$$\mathcal{L}_{\theta,\phi}(x) = \mathbb{E}_{q_\phi(z|x)}\left[\log p_\theta(x, z) - \log q_\phi(z|x)\right]$$
$$= \mathbb{E}_{p(\epsilon)}\left[\log p_\theta(x, z) - \log q_\phi(z|x)\right]$$

where $z = g(\epsilon, \phi, x)$.

- As a result we can form a simple Monte Carlo estimator of the individual-datapoint ELBO where we use a single noise sample $\epsilon$ from $p(\epsilon)$:

$$\epsilon \sim p(\epsilon)$$
$$z = g(\phi, x, \epsilon)$$
$$\tilde{\mathcal{L}}_{\theta,\phi}(x) = \log p_\theta(x, z) - \log q_\phi(z|x)$$

## Auto-Encoding Variational Bayes (AEVB) algorithm

**Algorithm 1:** Stochastic optimization of the ELBO. Since noise originates from both the minibatch sampling and sampling of $p(\epsilon)$, this is a doubly stochastic optimization procedure. We also refer to this procedure as the *Auto-Encoding Variational Bayes* (AEVB) algorithm.

---

**Data:**
    $\mathcal{D}$: Dataset
    $q_{\phi}(\mathbf{z}|\mathbf{x})$: Inference model
    $p_{\theta}(\mathbf{x}, \mathbf{z})$: Generative model

**Result:**
    $\theta, \phi$: Learned parameters

$(\theta, \phi) \leftarrow$ Initialize parameters
**while** *SGD not converged* **do**
    $\mathcal{M} \sim \mathcal{D}$ (Random minibatch of data)
    $\epsilon \sim p(\epsilon)$ (Random noise for every datapoint in $\mathcal{M}$)
    Compute $\tilde{\mathcal{L}}_{\theta, \phi}(\mathcal{M}, \epsilon)$ and its gradients $\nabla_{\theta, \phi} \tilde{\mathcal{L}}_{\theta, \phi}(\mathcal{M}, \epsilon)$
    Update $\theta$ and $\phi$ using SGD optimizer
**end**

---

## Computation of $log q_\phi(z|x)$

- Note that we typically know the density $p(\epsilon)$, since this is the density of the chosen noise distribution. As long as $g(.)$ is an invertible function, the densities of e and z are related as:

$$\log q_\phi(z|x) = \log p(\epsilon) - \log d_\phi(x, \epsilon)$$

where

$$\log d_\phi(x, \epsilon) = \log \left| \det \left( \frac{\partial z}{\partial \epsilon} \right) \right|$$

- The Jacobian matrix contains all first derivatives of the transformation from $\epsilon$ to $z$:

$$\frac{\partial z}{\partial \epsilon} = \frac{\partial(z_1, ..., z_k)}{\partial(\epsilon_1, ..., \epsilon_k)} = \begin{pmatrix} \frac{\partial z_1}{\partial \epsilon_1} & \cdots & \frac{\partial z_1}{\partial \epsilon_k} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_k}{\partial \epsilon_1} & \cdots & \frac{\partial z_k}{\partial \epsilon_k} \end{pmatrix}$$

- we can build very exible transformations $g()$ for which $\log d_\phi(\boldsymbol{x}, \epsilon)$ is simple to compute, resulting in highly exible inference models $q_\phi(\boldsymbol{z}|\boldsymbol{x})$.

## Factorized gaussian posteriors

A common choice is a simple factorized Gaussian encoder $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \mathrm{diag}(\sigma^2))$:

$$(\boldsymbol{\mu}, \log \sigma) = \mathrm{EncoderNeuralNet}_{\phi}(\mathbf{x})$$

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \prod_i q_{\phi}(z_i|\mathbf{x}) = \prod_i \mathcal{N}(z_i; \mu_i, \sigma_i^2)$$

where $\mathcal{N}(z_i; \mu_i, \sigma_i^2)$ is the PDF of the univariate Gaussian distribution.

- After reparameterization, we can write:

$$\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$$

$$(\boldsymbol{\mu}, \log \sigma) = \mathrm{EncoderNeuralNet}_{\phi}(\mathbf{x})$$

$$\mathbf{z} = \boldsymbol{\mu} + \sigma \odot \boldsymbol{\epsilon}$$

- The Jacobian of the transformation from $\epsilon$ to $z$ is:

$$\frac{\partial z}{\partial \epsilon} = \text{diag}(\sigma),$$

i.e. a diagonal matrix with the elements of $diag(\sigma)$ on the diagonal. The determinant of a diagonal (or more generally, triangular) matrix is the product of its diagonal terms. The log determinant of the Jacobian is therefore:

$$\log d_\phi(x, \epsilon) = \log \left| \det \left( \frac{\partial z}{\partial \epsilon} \right) \right| = \sum_i \log \sigma_i$$

and the posterior density is:

$$\begin{aligned} \log q_\phi(z|x) &= \log p(\epsilon) - \log d_\phi(x, \epsilon) \\ &= \sum_i \log \mathcal{N}(\epsilon_i; 0, 1) - \log \sigma_i \end{aligned}$$

when $z = g(\epsilon, \phi, x)$.

## Full-covariance Gaussian posterior

- The factorized Gaussian posterior can be extended to a Gaussian with full covariance:

$$q_\phi(z|x) = \mathcal{N}(z; \mu, \Sigma)$$

A reparameterization of this distribution is given by:

$$\epsilon \sim \mathcal{N}(0, I)$$
$$z = \mu + L\epsilon$$

- where **L** is a lower (or upper) triangular matrix, with non-zero entries on the diagonal.

The reason for this parameterization of the full-covariance Gaussian, is that the Jacobian determinant is remarkably simple. The Jacobian in this case is trivial: $\frac{\partial z}{\partial \epsilon} = L$. Note that the determinant of a triangular matrix is the product of its diagonal elements. Therefore, in this parameterization:

$$\log |\det(\frac{\partial z}{\partial \epsilon})| = \sum_i \log |L_{ii}|$$

And the log-density of the posterior is:

$$\log q_\phi(z|x) = \log p(\epsilon) - \sum_i \log |L_{ii}|$$

This parameterization corresponds to the Cholesky decomposition $\Sigma = LL^T$ of the covariance of $z$:

$$\Sigma = \mathbb{E}\left[(z - \mathbb{E}[z|])(z - \mathbb{E}[z|])^T\right]$$
$$= \mathbb{E}\left[L\epsilon(L\epsilon)^T\right] = L\mathbb{E}\left[\epsilon\epsilon^T\right]L^T$$
$$= LL^T$$

Note that $\mathbb{E}\left[\epsilon\epsilon^T\right] = I$ since $\epsilon \sim \mathcal{N}(0, I)$.

One way to build a matrix $\mathbf{L}$ with the desired properties, namely triangularity and non-zero diagonal entries, is by constructing it as follows:

$$(\boldsymbol{\mu}, \log \boldsymbol{\sigma}, \mathbf{L}') \leftarrow \text{EncoderNeuralNet}_{\phi}(\mathbf{x})$$
$$\mathbf{L} \leftarrow \mathbf{L}_{mask} \odot \mathbf{L}' + \text{diag}(\boldsymbol{\sigma})$$

$\mathbf{L}_{mask}$ is a masking matrix with zeros on and above the diagonal, and ones below the diagonal. The log-determinant is identical to the factorized Gaussian case:

$$\log \left| \det \left( \frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} \right) \right| = \sum_i \log \sigma_i$$

# VAE with a full-covariance Gaussian inference model

**Algorithm 2:** Computation of unbiased estimate of single-datapoint ELBO for example VAE with a full-covariance Gaussian inference model and a factorized Bernoulli generative model. $\mathbf{L}_{mask}$ is a masking matrix with zeros on and above the diagonal, and ones below the diagonal. Note that $\mathbf{L}$ must be a triangular matrix with positive entries on the diagonal.

**Data:**
    $x$: a datapoint, and optionally other conditioning information
    $\epsilon$: a random sample from $p(\epsilon) = \mathcal{N}(0, \mathbf{I})$
    $\theta$: Generative model parameters
    $\phi$: Inference model parameters
    $q_\phi(z|x)$: Inference model
    $p_\theta(x, z)$: Generative model
**Result:**
    $\mathcal{L}$: unbiased estimate of the single-datapoint ELBO $\mathcal{L}_{\theta,\phi}(x)$

$(\mu, \log\sigma, \mathbf{L}') \leftarrow \text{EncoderNeuralNet}_\phi(x)$
$\mathbf{L} \leftarrow \mathbf{L}_{mask} \odot \mathbf{L}' + \text{diag}(\sigma)$
$\epsilon \sim \mathcal{N}(0, \mathbf{I})$
$z \leftarrow \mathbf{L}\epsilon + \mu$
$\text{logqz} \leftarrow -\text{sum}(\frac{1}{2}(\epsilon^2 + \log(2\pi) + \log\sigma))$         $\triangleright = q_\phi(z|x)$
$\text{logpz} \leftarrow -\text{sum}(\frac{1}{2}(z^2 + \log(2\pi)))$         $\triangleright = p_\theta(z)$
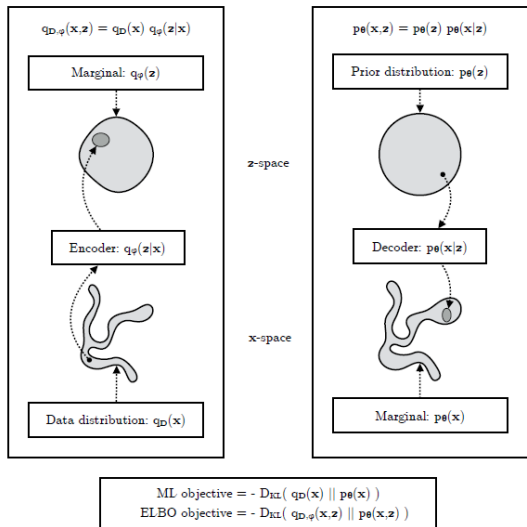$\mathbf{p} \leftarrow \text{DecoderNeuralNet}_\theta(z)$
$\text{logpx} \leftarrow \text{sum}(x \odot \log\mathbf{p} + (1-x) \odot \log(1-\mathbf{p}))$     $\triangleright = p_\theta(x|z)$
$\mathcal{L} = \text{logpx} + \text{logpz} - \text{logqz}$

## Estimation of the marginal likelihood

- After training a VAE, we can estimate the probability of data under the model using an importance sampling technique **Rezende et al. [2014]**

# Marginal likelihood and ELBO as KL divergences

## Marginal likelihood and ELBO as KL divergences

- One additional perspective is that the ELBOcan be viewed as a maximum likelihood objective in an augmented space.
- For some fixed choice of encoder $q_\phi(\boldsymbol{z}|\boldsymbol{x})$, we can view the joint distribution $p_\theta(\boldsymbol{x},\boldsymbol{z})$ as an augmented empirical distribution over the original data $\boldsymbol{x}$ and (stochastic) auxiliary features $\boldsymbol{z}$ associated with each datapoint.
- The model $p_\theta(\boldsymbol{x},\boldsymbol{z})$ then defines a joint model over the original data, and the auxiliary features.

## Other Examples of Variational treatment

- Variational Mixture of Gaussians
- Variational Linear Regression
- Variational Logistic Regression etc.
  **Bishop et al. [2006]**

**Thanks you.**
**Questions?**

## References

- Kingma, D.P. Variational Inference and Deep Learning: A New Synthesis (Ph.D. Thesis) (2017).

- Christopher M Bishop. Pattern recognition and machine learning, (2006).